

The Equidistance Index of Population Structure

Yaron Granot*¹, Saharon Rosset², and Karl Skorecki¹

¹Rappaport Faculty of Medicine and Research Institute, Technion–Israel Institute of Technology, and Rambam Medical Center, Haifa, Israel

²School of Mathematical Sciences Tel Aviv University, Tel Aviv, Israel

*E-mail: yarongranot@hotmail.com

Abstract

Measures of population differentiation, such as F_{ST} , are traditionally derived from the ratio of genetic diversity within and between populations. However, the emergence of population clusters from multilocus analysis is a function of genetic *structure* (departures from panmixia) rather than diversity. If the populations are close to panmixia, slight differences between the mean pairwise distance within and between populations (low F_{ST}) can manifest as strong separation between the populations, thus population clusters are often evident even when the vast majority of diversity is partitioned within populations rather than between them. Moreover, because F_{ST} utilizes the mean rather than *deviations* from the mean, it does not directly reflect the strength of separation between population clusters. For any given F_{ST} value, clusters can be tighter (more panmictic) or looser (more stratified), and in this respect higher F_{ST} does not always imply stronger differentiation. In this study we propose substituting the mean in the F_{ST} equation with the standard deviation, thereby deriving a novel measure of population separability, denoted E_{ST} , which is more consistent with clustering and classification. To assess the utility of this metric, we ranked various human (HGDP) population pairs based on F_{ST} and E_{ST} and found substantial differences in ranking order. In some cases examined, most notably among isolated Amazonian tribes, E_{ST} ranking seems more consistent with demographic, phylogeographic and linguistic measures of classification compared to F_{ST} . Thus, E_{ST} may at times outperform F_{ST} in identifying evolutionarily significant differentiation.

Keywords: F_{ST} , population structure, panmixia, differentiation, standard deviation, genetic isolates

Introduction

Genetic differentiation among populations is typically derived from the ratio of within- to between-population diversity. The most commonly used metric, F_{ST} , was originally introduced as a fixation index at a single biallelic locus (Wright 1978), and subsequently adapted as a measure of population subdivision by averaging the values over many loci (Nei 1973; Weir and Cockerham 1984). F_{ST} can be expressed mathematically as $F_{ST}=1-S/T$, where S and T represent heterozygosity or some other measure of diversity in subpopulations and in the total population (Hudson et al. 1992). The validity of F_{ST} as a measure of differentiation has been brought into question, especially when gene diversity is high (e.g., in microsatellites), and various metrics, including G'_{ST} (Hedrick 2005) and Jost's D (Jost 2008), have been proposed to address this inadequacy (though see Whitlock 2011 for a counter-perspective).

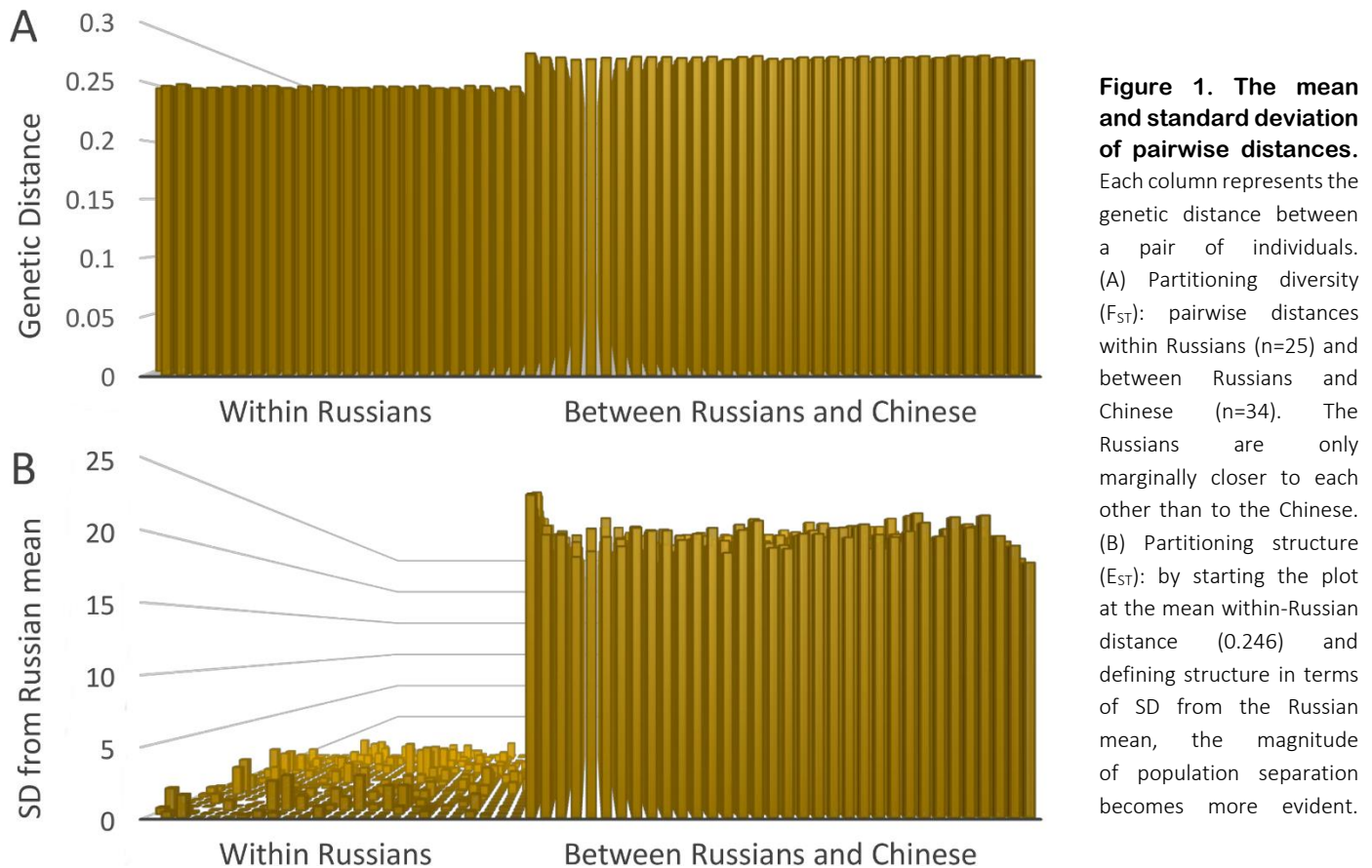
Although these metrics vary considerably in their formulation, they all follow the same basic framework of partitioning genetic diversity into within- vs. between-group components. It has long been noted, however, that the apportionment of diversity (Lewontin 1972) does not directly reflect the strength of separation between populations, and the emergence of population clusters has been shown both empirically (Mitton 1977) and mathematically (Edwards 2003; Tal 2013) even when the vast majority of diversity is within rather than between populations. For example, humans sampled from across Europe (Nelis et al. 2009) and East Asia (Tian et al. 2008) form identifiable clusters with pairwise F_{ST} as low as 0.002, even though 99.8% of the variation is contained within populations and only 0.2% is between them. Clearly, these clusters reflect an aspect of population differentiation that is not directly captured by F_{ST} , yet there is currently no commonly used metric for partitioning structure into within- and between-population components in the same way that F_{ST} partitions diversity. Clustering algorithms such as principal component analysis (PCA) (Patterson et al. 2006) and STRUCTURE (Rosenberg et al. 2002) are widely circulated, however such programs are primarily used for visualization, and there is still value in summary statistics for quantifying complex datasets on a simple 0-1 scale.

Here we propose a novel statistic, denoted E_{ST} , based on a modified F_{ST} estimator in which the mean pairwise distance between individuals (a measure of diversity) is replaced by the standard deviation of pairwise distances (a measure of structure), thus extracting the excess structure in the total population compared to subpopulations. Conceptually, E_{ST} is formulated in three steps: 1. Population structure is defined in terms of departures from *panmixia*. 2. Panmixia is defined in terms of pairwise *equidistance* between individuals (a population is considered panmictic if all individuals are equally distant from each other). 3. Departures from equidistance are defined in terms of the *standard deviation* of pairwise distances. E_{ST} reflects the decrease in panmixia when subpopulations are pooled. The general formula is: $E_{ST} = 1 - SD_S / SD_T$, where SD_S and SD_T represent the standard deviations of pairwise distances in subpopulations and in the total population. While F_{ST} is weighed down by high *diversity* within populations, E_{ST} is weighed down by high *structure* within populations. Since diversity is usually greater than structure, E_{ST} is usually greater than F_{ST} .

Results and Discussion

Partitioning Diversity vs. Partitioning Structure

The difference between the partitions of diversity and structure within and between two populations from the human genome diversity project (HGDP) (Cann et al. 2002) is illustrated in Figure 1. The mean distance among mixed Russian-Chinese pairs is only marginally (~10%) higher than among Russian-Russian pairs (Figure 1A), reflecting the relatively low F_{ST} . However this translates to a far greater increase in total structure compared to the low structure within each population (Figure 1B), reflecting the much higher E_{ST} .



We compared F_{ST} , E_{ST} , and clustering among Russian and Chinese samples, with an increasing amount of single nucleotide polymorphisms (SNPs) ranging from 10 to 660,755 (Figure 2). Using multidimensional scaling (MDS), the two population clusters gradually diverge as SNP count increases, with no corresponding increase in F_{ST} . At the same time we observe a steady increase in E_{ST} directly corresponding to the emerging clusters, indicating that the Russian and Chinese HGDP samples are close to panmixia. With few SNPs this is obfuscated by the variance of the genetic distance measure, hence E_{ST} is relatively small. The actual levels of panmixia become increasingly evident as more SNPs are added, thus revealing the population clusters (Edwards 2003). However this process does not proceed indefinitely; the finite number of pairwise differences among humans (~3 million SNPs) sets an upper limit to the number of available markers, and the amount of extractable information is further reduced by linkage disequilibrium. In our data the increase in E_{ST} as a function of marker count reaches a plateau above 100,000 SNPs (Figure S1). Although this upper bound can vary across different datasets and types of markers, it suggests that resolution may not improve substantially with further increases in marker count. Thus, these clusters can be considered close approximations of the “true” strength of separation among these populations. For this reason, E_{ST} estimates should include as many markers as possible, although fewer markers can be used and the terminal E_{ST} can be extrapolated.

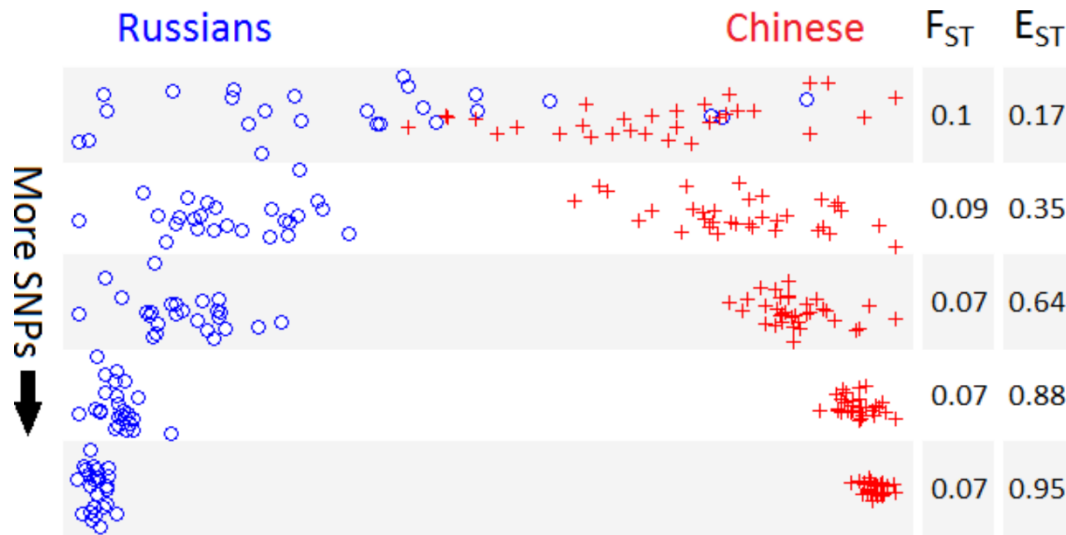


Figure 2. F_{ST} and E_{ST} vs. Clustering with increasing SNP count.

Multidimensional scaling (MDS) plots with Russian ($n=25$) and Chinese ($n=34$) samples with increasing SNP count from top to bottom (10, 100, 1000, 10,000, and 660,755 SNPs). Two clusters gradually emerge as SNP count increases, along with an increase in E_{ST} , while F_{ST} remains relatively constant.

In order to determine whether or not E_{ST} adds insight to the analysis of population structure, we sought to compare the rank order of population differentiation using F_{ST} and E_{ST} . Pairwise F_{ST} and E_{ST} values from various HGDP populations are given in Table 1 (see Table S1 and Figure S2 for additional comparisons). As expected, $E_{ST} > F_{ST}$ in most population pairs. Only the Colombian-Maya pair has a slightly lower E_{ST} than F_{ST} , due to a combination of relatively low differentiation and high levels of intra-population structure. According to the HGDP browser (http://spsmart.cesga.es/search.php?dataSet=ceph_stanford, the Colombians ($n=7$) are the only HGDP population sample where two different tribes (Piapoco and Curripaco) were combined, which can help explain the high level of structure observed in this particular population (see Table S1, Figure S10, and Materials and Methods for further analysis of E_{ST} range).

Table 1

Pairwise F_{ST} (above diagonal) and E_{ST} (below diagonal) in 5 New World and 5 Old World HGDP populations

	Surui	Karitiana	Colombian	Maya	Pima	Yakut	Mongola	Russian	Bantu	San
Surui		0.13	0.1	0.09	0.12	0.15	0.15	0.17	0.23	0.3
Karitiana	0.58		0.08	0.07	0.11	0.13	0.13	0.16	0.22	0.29
Colombian	0.51	0.57		0.03	0.06	0.09	0.09	0.12	0.18	0.25
Maya	0.52	0.63	0.02		0.04	0.07	0.06	0.08	0.15	0.21
Pima	0.57	0.63	0.37	0.43		0.1	0.09	0.12	0.19	0.25
Yakut	0.74	0.8	0.6	0.69	0.74		0.01	0.06	0.13	0.19
Mongola	0.81	0.87	0.69	0.8	0.83	0.46		0.06	0.12	0.19
Russian	0.82	0.87	0.74	0.83	0.84	0.86	0.9		0.11	0.17
Bantu	0.88	0.92	0.85	0.91	0.91	0.93	0.94	0.95		0.07
San	0.92	0.95	0.89	0.95	0.94	0.96	0.97	0.98	0.89	

Amazonians vs. Global Populations

The Surui and Karitiana have an unusually high pairwise F_{ST} . In fact, the Karitiana are as diverged from the neighboring Surui in terms of F_{ST} as they are from the Mongola on the other side of the world (Table 1, Figure 3, and Figure S6). Moreover, F_{ST} actually decreases initially with distance from the Amazon, from 0.13 between the two Amazonian tribes, to 0.08-0.1 between Amazonians and Colombians, further decreasing to 0.07-0.09 between Amazonians and the more distant Maya. Remarkably, the highest F_{ST} among all HGDP Native American populations is between the two geographically closest populations, the Surui and Karitiana. These apparent anomalies can be explained by the inflation of F_{ST} in genetic isolates. F_{ST} between pairs of isolates can be nearly twice as high as between either one of the isolates and a more cosmopolitan population, as pairwise F_{ST} reflects the *combined* isolation of both populations. Since the Surui and Karitiana are both isolated, their pairwise F_{ST} is nearly double that between any one of them and a larger, less isolated population such as the Maya. In other words, the Maya's contribution to the pairwise F_{ST} is dwarfed by that of the Amazonians.

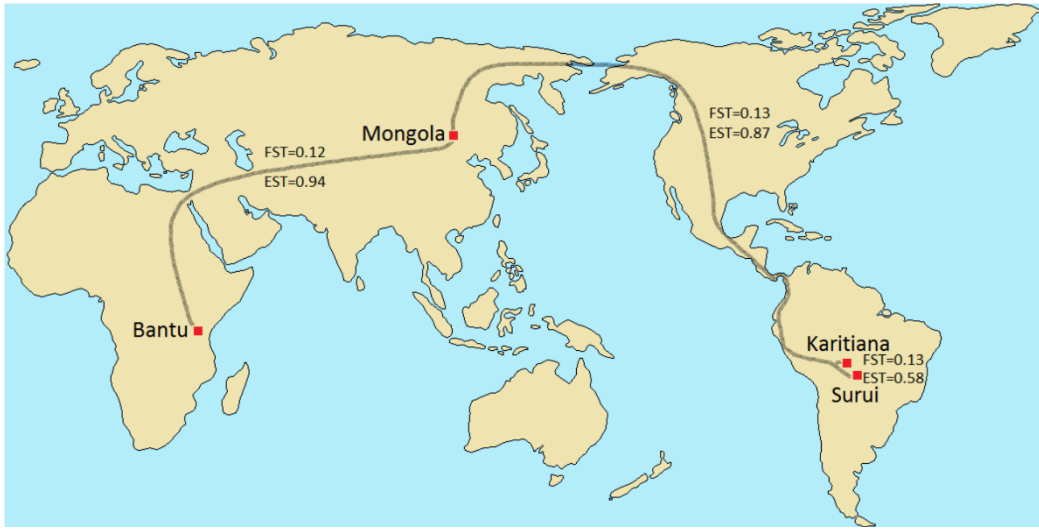


Figure 3. Geographic distance vs. F_{ST} and E_{ST} in various populations.

In terms of F_{ST} , the Karitiana are roughly as diverged from the nearby Surui ($F_{ST}=0.13$) as they are from the Mongola on the other side of the world ($F_{ST}=0.13$) or as the Bantu are from the Mongola ($F_{ST}=0.12$). In terms of E_{ST} , differentiation is far greater among these global populations ($E_{ST}\approx 0.9$) than between the neighboring Amazonian tribes ($E_{ST}\approx 0.6$).

Differentiation based on E_{ST} (Surui-Karitiana=0.58, Karitiana-Mongola=0.87, and Mongola-Bantu=0.94) seems more consistent with the geographic distances among these populations (Figure 3). It should be noted that the Surui-Karitiana E_{ST} might be somewhat underestimated due to cryptic sampling of close relatives (Rosenberg 2006), however the wide range of heterozygosity values (which are less sensitive to the sampling of close relatives) and the elevated structure across all Native American HGDP populations (Figures S3-S5) suggest that this is not merely a sampling artifact. In some cases E_{ST} also decreases with distance from the Amazon (Table 1), however this decrease is more moderate than the decrease in F_{ST} (Figure S6).

Neighbor-joining trees of individual similarities (Jorde and Wooding 2004) are a convenient tool for representing multidimensional genetic data on a two-dimensional plane, while simultaneously displaying distances within and between populations. Two pairs of such trees, for Surui-Karitiana and Yoruba-Russians, are given in Figure 4, and we can see that in both cases distances are greater between individuals (black branches) than between populations (red branches) (Figure 4A).

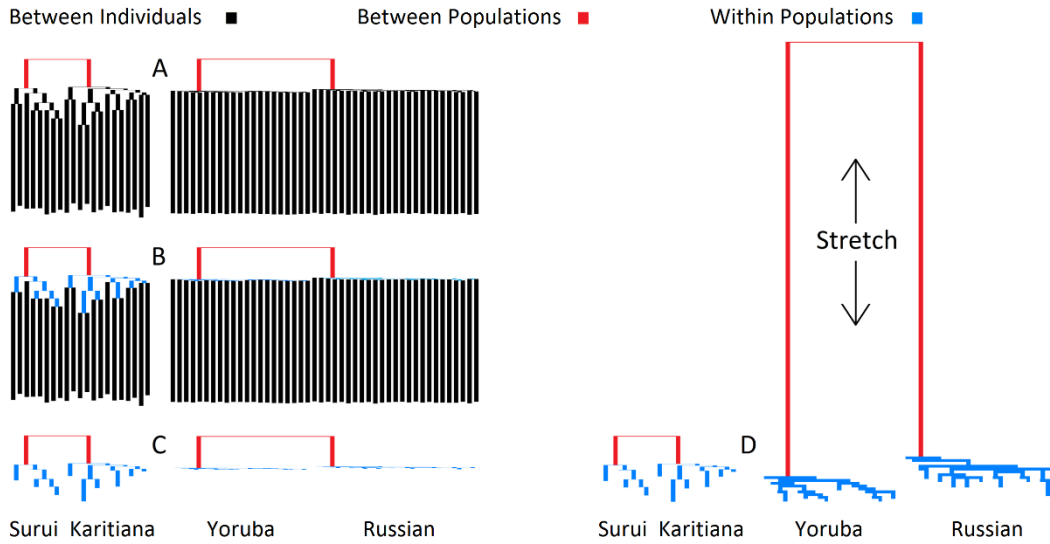


Figure 4. Surui-Karitiana vs. Yoruba-Russian NJ trees of individual similarities.

(A) Diversity is apportioned into individual (black) and population (red) components. (B) A third component, structure within populations (blue), is added. (C) The individual component is removed. (D) The Yoruba-Russian tree is stretched to roughly match the level of structure within the Surui-Karitiana tree.

The ratio of within- to between-population distance is roughly equivalent in the two population pairs, however the Yoruba-Russian tree is significantly *flatter*, indicating greater panmixia within these two populations (Figures S7-S8). Adding a third dimension of intra-population structure (blue branches) highlights this discrepancy (Figure 4B), which is further accentuated by removing the inter-individual component (Figure 4C) and stretching the Yoruba-Russian tree to match the level of structure observed in the Surui-Karitiana tree (Figure 4D). At first glance the Amazonian tribes, with their long population branches, appear to be as differentiated as the Yoruba are from the Russians. Upon closer inspection, however, the Yoruba and Russians appear more strongly diverged. The Amazonian tribes are highly structured not only between them, but also within them, resulting in distant, but loosely separated clusters. This aspect of population structure is not captured by F_{ST} , which is actually slightly higher between the Surui and Karitiana (0.13) than between Yoruba and Russians (0.12), but is revealed by the higher E_{ST} between Yoruba and Russians (0.97) compared to the Surui and Karitiana (0.58).

E_{ST} and the Dissimilarity Fraction

The dissimilarity fraction, ω , is defined (Witherspoon et al. 2007) as the probability that individuals are genetically more similar to members of a different population than to members of their own population. For pairs of populations, this probability should have a 0-0.5 range, with $\omega=0$ indicating that individuals are always closer to members of their own population and $\omega=0.5$ indicating that individuals are just as likely to be closer to members of the other population as to members of their own population. Witherspoon et al. reported that that when many thousands of loci are analyzed, individuals from “geographically separated populations” are never closer to each other than to members of their own populations. The definition of “geographically separated” is, of course, open to interpretation. We found no overlap ($\omega=0$) between the Adygei and Uygur HGDP samples, but some overlap ($\omega > 0$) between Mayans and Surui, despite a 4x higher F_{ST} (Figure 5). Thus, F_{ST} and the dissimilarity fraction (ω) are not necessarily congruent. The E_{ST} values for these two population pairs are more consistent with ω , showing strong separation between the Adygei and Uygur (0.79) and more moderate separation between Colombians and Maya (0.52) (see Figure S9 for a more detailed plot).

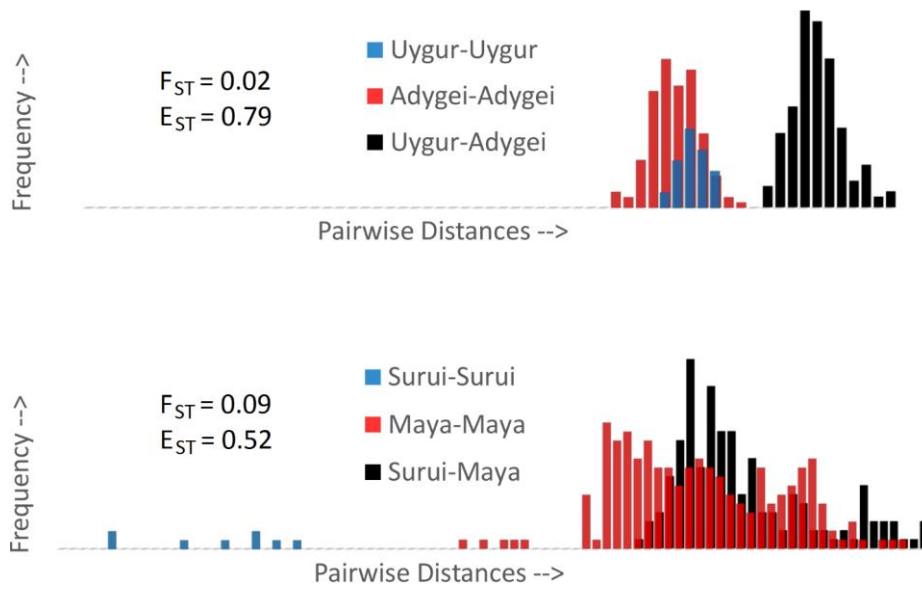


Figure 5. F_{ST} vs. genetic similarity in various population pairs.

Pairwise distances are colored red or blue within populations and black between populations. (A) Even at a relatively low F_{ST} of 0.02 all within-population pairs among the Uygur and Adygei samples are genetically more similar than all the between-population pairs. (B) Separation is more ambiguous among Native Americans. Despite a relatively high F_{ST} of 0.09, there is substantial overlap between Maya-Maya (red) and Maya-Surui (black) samples. E_{ST} values are more consistent with the within- vs.-between population overlap and the dissimilarity fraction (w).

Summary and Conclusions

The main distinction between F_{ST} and E_{ST} is that F_{ST} partitions diversity, whereas E_{ST} partitions structure within and between populations. F_{ST} is more sensitive to effective population size, while E_{ST} is more sensitive to outliers, though this is largely mitigated by using E_{ST} median rather than E_{ST} mean (see Materials and Methods). F_{ST} is often weighed down by high levels of intrapopulation diversity and can be close to zero even when population clusters are completely separated. This is not a flaw in F_{ST} , but it does demonstrate a conceptual disconnect between F_{ST} and clustering. Sewall Wright proposed a series of arbitrary F_{ST} thresholds ranging from 0.05 to 0.25, denoting little to very great differentiation (Wright 1978), however these are only broad guidelines, and the highest ranking of “very great differentiation” leaves most of the range (0.25-1) undefined.

Given its wider empirical range and more direct correlation with clustering and classification (Figure 2), phylogeography (Figure 3), and the dissimilarity fraction (Figure 5), such arbitrary thresholds may not be necessary for E_{ST} . $E_{ST} > 0.5$ simply indicates that most of the structure is between populations rather than within them, corresponding to moderately separated populations such as Russians and Adygei ($E_{ST}=0.5$), Bantu from South Africa and Kenya ($E_{ST}=0.48$), or French and Sardinians ($E_{ST}=0.48$) (Table S1). $E_{ST} < 0.5$ indicates weak differentiation and $E_{ST} >> 0.5$ indicates strong differentiation. E_{BT} is similar in many ways to E_{ST} , though its HGDP ranking order is often intermediate between F_{ST} and E_{ST} (Table S1). Interestingly, some East Asians populations have relatively low E_{BT} , such as Cambodians vs. Mongola ($E_{BT}=0.13$) and Japanese vs. Chinese ($E_{BT}=0.16$).

Differentiation metrics are judged by their ability to quantify meaningful evolutionary divergence, and can be indispensable in identifying *Evolutionarily Significant Units* (ESU) and *Distinct Population Segments* (DPS) for conservation (Waples 1991). For example given several subpopulations within a species, it is reasonable to prioritize the most highly differentiated subpopulation for conservation in order to maximize biodiversity. However, higher F_{ST} does not necessarily reflect stronger separation and lower misclassification, as with the Uygur and Adygei, whose clusters are better defined than those of the Surui and Maya despite a fourfold lower F_{ST} (Figure 5). In this

context humans can be a useful model species simply because we know so much about human populations due to our “long habit of observing ourselves” (Darwin 1871). This allows us to make educated inferences about human populations that might otherwise be overlooked, e.g., we can be skeptical of the high Surui-Karitiana F_{ST} , and realize that this is most likely due to the relatively recent isolation of two small tribes. This is a luxury that we do not usually have with other species, in which case high F_{ST} can be misinterpreted as a deep phylogenetic divide, potentially leading to misguided conservation strategies. Our hope is that by combining information from both *fixation* (F_{ST}) and *equidistance* (E_{ST}) indices, researchers could make more informed decisions.

Unlike F_{ST} , which can be estimated from a handful of markers, E_{ST} requires large datasets with thousands of markers, which were unavailable to previous generations of population geneticists. With the latest SNP chips containing well over 100,000 markers, accurate estimates of departures from panmixia are finally within reach, and there is no longer a need for the simplifying assumption that subpopulations are effectively panmictic. By deriving an F_{ST} -type statistic for apportioning structure within and between populations, namely E_{ST} , we hope to add a new useful metric to the 21st century population genetics toolkit.

Materials and Methods

The HGDP data used in our analysis are available at: <http://www.hagsc.org/hgdp/files.html>. After removing the 163 mitochondrial SNPs and 105 samples previously inferred to be close relatives (Rosenberg 2006), the final file included 660,755 SNPs from 938 samples in 53 populations. Strings of SNPs were treated as sequences, with mismatches summed and divided by the sequence length. Pairwise distances, based on Allele Sharing Distance (ASD) (Gao and Martin 2009), were calculated as one minus half the average number of shared alleles per locus.

We used Hudson’s F_{ST} estimator (Hudson et al. 1992):

$$F_{ST} = 1 - S/T \quad (1)$$

Where S and T are the mean pairwise distances within subpopulations and in the total pooled population.

The general equation for E_{ST} is:

$$E_{ST} = 1 - SD_S / SD_T \quad (2)$$

Where SD_S and SD_T are the standard deviations (SD) of pairwise distances within subpopulations and in the total population. This E_{ST} estimator is referred to as $E_{ST}mean$. We used three additional E_{ST} estimators: $E_{ST}min$, $E_{ST}median$, and $E_{ST}max$ (Figure S10). All four estimators use the same basic formula, with only the type of SD_S differing among estimators. In $E_{ST}min$, $E_{ST}median$, and $E_{ST}max$, SD_S is respectively replaced with the smallest, median, and largest *individual SD*, where the individual SD is the standard deviation of pairwise distances between a single sample and all other samples in the population. $E_{ST}min$ uses the smallest individual SD_S from each population, i.e., the SD of the most panmictic sample, $E_{ST}median$ uses the median individual SD_S , and $E_{ST}max$ uses the highest individual SD_S . Each of these

metrics has different sensitivities to various sampling biases. Due to E_{STmean} 's sensitivity to the sampling of close relatives, we used $E_{STmedian}$ (which is unaffected by the inclusion of relatives as long as at least 50% of the samples are unrelated) as the primary measure of E_{ST} in this study. In the rare event that >50% of the samples are closely related, E_{STmax} may be preferable, as long as at least one individual has no close relatives among the samples. E_{ST} values, especially E_{STmin} and E_{STmean} , can be negative if structure is high and differentiation is low (Figure S10). Small sample sizes were often sufficient for estimating heterozygosity (Figure S11) and F_{ST} and E_{ST} (Figure S12) using all the SNPs in the HGDP dataset.

We derived an additional equidistance index, denoted E_{BT} , which is less sensitive to intra-population structure and the inclusion of relatives. Recall that E_{ST} reflects equidistance (E) within subpopulations (S) compared to the total (T) population. Similarly, E_{BT} reflects equidistance (E) between subpopulations (B) compared to the total (T) population:

$$E_{BT} = 1 - SD_B / SD_T \quad (3)$$

Where SD_B and SD_T are the standard deviations of pairwise distances between individuals from different subpopulations, and in the total pooled population. In most cases $SD_T \geq SD_B$, because SD_T includes pairs of individuals from the same population as well as pairs from different populations, whereas SD_B only includes pairs of individuals from different populations. Pairs of individuals from the same population are likely to have a higher SD due to relatives in the samples, which disrupt the panmixia (see Naxi population in Figures S3-S5). Panmictic populations are not just equidistant among themselves, they are also equidistant towards each other. Such populations should have similar SD_S and SD_B , and thus similar E_{ST} and E_{BT} . All F_{ST} , E_{ST} and E_{BT} estimates in this study are based on pairwise comparisons between two populations or population groups. Each of the two paired populations was given equal weight, as were the within- and between-population pairs. Thus, 25% of the total weight was given to each population, and 50% to between-population pairs.

We developed a custom MATLAB code for extracting genetic distances from SNP data and estimating heterozygosity, pairwise distances, F_{ST} , E_{ST} , and E_{BT} . The code corrects for missing data and small sample sizes, and identifies outliers, but includes no further assumptions or corrections. Phylogenetic trees and MDS plots were also generated with MATLAB. Equal angle and square neighbor-joining trees of individual similarities were generated from matrices of pairwise distances with the *seqneighjoin* command. An alternative script, based on the internal MATLAB *seqpdist* command for sequence distance, yielded similar results.

Acknowledgments

We thank Alan Templeton for helpful advice, Omri Tal for mathematical input, and Sagi Abelson for help with the MATLAB script.

References

1. Wright S (1978) Evolution and the Genetics of Populations. Vol. 4, Variability Within and Among Natural Populations. University of Chicago Press, Chicago.
2. Nei M (1973) Analysis of gene diversity in subdivided populations. Proc Natl Acad Sci USA, 70, 3321–3323.

3. Weir BS, Cockerham CC (1984) Estimating F-Statistics for the Analysis of Population Structure. *Evolution* 38 (6): 1358.
4. Hudson RR, Slatkin M, Maddison WP (1992) Estimation of levels of gene flow from DNA sequence data. *Genetics* 132 (2), 583–589.
5. Hedrick PW (2005) A standardized genetic differentiation measure. *Evolution*, 59, 1633–1638.
6. Jost L (2008) G_{ST} and its relatives do not measure differentiation. *Mol Ecol*, 17, 4015–4026.
7. Whitlock M (2011) G'_{st} and D do not replace F_{st} . *Mol Ecol* 20:1083–109.
8. Lewontin RC (1972) The apportionment of human diversity. *Evol Biol*, 6, 381–98.
9. Mitton JB (1977) Genetic Differentiation of Races of Man as Judged by Single-Locus and Multilocus Analyses. *Amer Nat*, 111 (978), 203–212.
10. Edwards AWF (2003) Human genetic diversity: Lewontin's fallacy, *Bioessays*, 25, 798–801.
11. Tal O (2013) Two complementary perspectives on inter-individual genetic distance. *Biosystems*, 111: 18–36.
12. Nelis M, Esko T, Mägi R, Zimprich F, Zimprich A, Toncheva D, et al. (2009) Genetic structure of Europeans: a view from the North-East. *PLoS One* 4: e5472.
13. Tian C, Kosoy R, Lee A, Ransom M, Belmont JW, Gregersen PK, et al. (2008) Analysis of East Asia Genetic Substructure Using Genome-Wide SNP Arrays. *PLoS ONE*, 3, (12): e3862.
14. Patterson N, Price AL, Reich D (2006) Population Structure and Eigenanalysis. *PLoS Genet* 2: 2074–2093. doi:10.1371/journal.pgen.0020190.
15. Rosenberg NA, Pritchard JK, Weber JL, Cann HM, Kidd KK, e Zhivotovsky LA, et al. (2002) Genetic structure of human populations. *Science* 298: 2381–2385.
16. Cann HM, de Toma C, Cazes L, Legrand MF, Morel V, Piouffre L (2002) A human genome diversity cell line panel. *Science* 296 (5566): 261–2.
17. Rosenberg NA (2006) Standardized subsets of the HGDP-CEPH Human Genome Diversity Cell Line Panel, accounting for atypical and duplicated samples and pairs of close relatives. *Ann Hum Genet*, 70: 841–847.
18. Jorde LB, Wooding SP (2004) Genetic variation, classification, and "race." *Nat Genet*, 36, S28–32.
19. Witherspoon DJ, Wooding S, Rogers AR, Marchani EE, Watkins WS, Batzer MA, et al. (2007) Genetic Similarities Within and Between Human Populations. *Genetics* 176(1): 351–9. pmid:17339205 doi: 10.1534/genetics.106.067355.
20. Waples RS (1991) Definition of 'species' under the Endangered Species Act: Application to Pacific salmon. U.S. Department of Commerce NOAA Technical Memorandum, NMFS, F/NWC–194.
21. Darwin C (1871) *The Descent of Man and Selection in Relation to Sex*. London: John Murray.
22. Gao X, Martin ER (2009) Using allele sharing distance for detecting human population stratification. *Hum Hered*, 68, 182–191.